

Identification of depression and PTSD among Twitter users using pre-trained language model

Anonymous EMNLP submission

Abstract

001 Suicide is a global health issue and early diagnosis is necessary for effective treatment. Recent advancements in natural language processing has aided the identification of mental health disorders in social media. This paper investigated the efficacy of pre-trained language model (PLM) in identifying depression and post-traumatic stress disorder (PTSD) with Twitter data. Leveraging the CLPsych 2015 dataset (which constitutes of tweets from users with depression, PTSD and neither condition), we implemented various experimental designs using Long Short Term Memory (LSTM) and attention. The results demonstrate that while performance decreases for multi-nominal classification, the detection of mental health conditions improves with the implementation of attention. This study also underscores the complexity of differentiating between overlapping lexicons with multiple mental health conditions and highlights the potential of PLMs in supporting mental health diagnosis.

023 1 Introduction

024 Suicide is a global health problem and is the fourth leading cause of death for the 15-44 years demographic globally (World Health Organization, 2021). Mental disorders, including depression and post-traumatic stress disorder (PTSD) have been found to increase the likelihood of suicidal ideation and suicide (Holliday et al., 2021; Busby Grant et al., 2023; Chou et al., 2023; Kratovic et al., 2021). These disorders not only hamper the quality of life for the people who suffer with them but also lessen the quality of life for their families and environment (García-Noguez et al., 2023). Moreover, 75% of people with a severe mental disorder do not receive treatment (Ji et al., 2021). Early diagnosis and subsequent treatment can help to lessen the negative impacts that arise from mental health disorders (Beirão et al., 2020; Kearns et al., 2012).

041 Researchers are leveraging social context to better understand mental health problems and has been an ongoing process. In the past, researchers used Google trends for mental health surveillance (Page et al., 2011), examining depression based chatter on Twitter (Cavazos-Rehg et al., 2016) and implementing machine learning algorithms to classify tweets in terms of stress or relaxation (Doan et al., 2017). Recently, advancements in natural language processing (NLP) and pre-trained language models (PLMs) have been helpful in identifying the mental health disorder traits from textual data (Ji et al., 2021; Vajre et al., 2021). Although these methods will never fully replace the psychiatric diagnosis and psychotherapy, they assist researchers and clinicians in early detection of mental health symptoms.

058 Prior to the advancement of PLMs, an early study was conducted in 2014 as a part of a hackathon event (Coppersmith et al., 2014). The authors performed a binary classification between the combinations of control, PTSD and depression outcomes based on the tweets gathered via Twitter api (Coppersmith et al., 2015). Following this research, the same dataset has aided other research, for example, interpreting mental health outcomes (Yang et al., 2023), training new PLMs centric to mental health outcomes (Ji et al., 2021) and comparing various machine learning models for their effectiveness in capturing mental health outcomes (Husseini Orabi et al., 2018).

072 However, the aforementioned studies focused on binary classification (depression vs control group) to identify the presence or absence of depression among Twitter users. Although there are overlapping expressions between PTSD and depression, there are also dissimilarities between the two mental disorders. Given how these disorders may affect an individual differently, identification of PTSD and depression separately could influence an individual's journey to recovery (Finch, 2023). Proper

082 diagnosis allows clinicians to recommend thera- 130
083 peutic interventions based on specific conditions 131
084 (Finch, 2023; Kimberly Holland, Timothy J. Legg, 132
085 2019). As such, in this research, we extend the 133
086 classification to all categories of CLPsych 2015 134
087 dataset, i.e. depression, PTSD and control, based 135
088 on tweets. 136

089 2 Methodology 137

090 We aim to answer two key questions in this paper: 138
091 1. How effective are PLMs for tracking multiple 139
092 mental health problems? 2. Which method is most 140
093 effective for handling multi-nominal mental health 141
094 classification? 142

095 Alongside the two questions, we also scrutinise the 143
096 scenarios where only depression detection or the 144
097 detection of general mental health issues might be 145
098 essential. 146

099 2.1 CLPsych 2015 shared dataset 147

100 The CLPsych 2015 shared dataset contains publicly 148
101 available tweets collected from the Twitter api over 149
102 the period 2008 to 2013. The tweets were posted by 150
103 users with PTSD, depression and a control group 151
104 who did not have any identified mental health con- 152
105 ditions as per tweets (Coppersmith et al., 2014). In 153
106 total, there are 1145 training set and 599 testing set 154
107 of anonymous users. Please note that the numbers 155
108 may not match the original set due to the exclusion 156
109 of users whose conditions were not recorded. 157
110 For this study, we used all available users and their 158
111 subsequent tweets to identify their category of men- 159
112 tal health condition, if present. Since the number 160
113 of control (572 training, 299 testing) users were 161
114 higher than depression (327 training, 150 testing) 162
115 and PTSD (246 training, 150 testing) users, we 163
116 used weighted cross entropy function for calcula- 164
117 tion of loss. However, the number of tweets was 165
118 reduced to a maximum of last 1000 tweets per user 166
119 out of a possible maximum of 3000 tweets per user 167
120 due to computational constraints. Despite this, each 168
121 epoch per experiment took over a day due to the 169
122 large volume of the dataset and the reliance on the 170
123 sequential computation of the tweets.

124 2.2 Algorithm for the experimental designs 171

125 The experiments were run for all users using Algo- 172
126 rithm 1. The number of epochs was set to 20, with 173
127 the training loop exiting if there was a increase 174
128 in the training loss. A single user was taken as 175
129 their own batch for training because of the choice

of model designs. Please refer to Section 3 for 130
the model designs. All the tweets went through 131
pre-processing phase where the textual content was 132
cleaned removing any white spaces, retweets, men- 133
tions, URLs, punctuation and emoticons. Please 134
note that cross-validation was not feasible due to 135
the magnitude of the dataset. 136

For each user u_i , their individual tweets 137
 t_1, t_2, \dots, t_n were tokenized and passed through 138
a pre-trained RoBERTa model. The details of the 139
choice of PLM is provided in section 2.3. The out- 140
put was a tensor containing the embedding of the 141
tweet t_i . The 768 dimension $[CLS]$ token, which 142
contains the classification information of the en- 143
tire sentence (Devlin et al., 2018), was extracted 144
for each tweet. For each user, these $[CLS]$ to- 145
kens were then stacked to form the tensor of shape 146
 $t_n^{u_i} \times 768$, where $t_n^{u_i}$ was the number of tweets 147
for user u_i . Further experiments were performed 148
using these stacked tensors as explained in Section 149
3. The output of each experiment was then con- 150
nected to two fully connected layers, with $\tanh()$ 151
as the activation function on both layers. The first 152
layer converted the output from 768 dimensions 153
to 100 dimensions and the second layer converted 154
from 100 dimensions to 3 dimensions. The output 155
of the second fully connected layer was passed to 156
softmax function, given by, $\sigma(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$, to 157
convert the results into probabilities. The final out- 158
put was the category (control, depression or PTSD) 159
with the highest probability i.e. $\max(\sigma(x_i))$. 160

Algorithm 1 Training CLPsych 2015 dataset

```
for epochs ( $e_i$ ) = 1 to  $e$  do
  for users ( $u_i$ ) = 1 to  $u$  do
    Pre-process each tweet removing any punctu-
    ation, white space, links, retweets and
    emoticons
    Pass tweet to tokenizer and pre-trained
    RoBERTa and extract  $[CLS]$  token
    Stack all  $[CLS]$  tokens for user  $u_i$ 
    Perform experiment  $E$  on stacked  $[CLS]$ 
    embedding
    Two layers of MLP with  $\tanh()$  and
     $\text{softmax}()$  to compute predicted  $\hat{y}$ 
    Calculate loss and update weight
  end for
end for
```

161	2.3 RoBERTa for base embeddings	209
162	We used a Twitter-based fine-tuned model of	210
163	RoBERTa called <i>cardiffnlp/twitter-roberta-base</i>	211
164	(Barbieri et al., 2020) for the base embeddings as	212
165	our PLM. The embeddings were extracted using	213
166	<i>transformer</i> library (Wolf et al., 2019). The smaller	214
167	memory size of RoBERTa and its pre-training on	215
168	Twitter data made it an appropriate choice for this	216
169	study. There was an expectancy that the localisa-	217
170	tion of Twitter vocabulary was present in the PLM	218
171	of choice. Therefore, it provided appropriate token	
172	embeddings for further experiments.	
173	3 Experimental Designs	219
174	We describe four implemented network models	220
175	which were used to evaluate the performance of	221
176	the detection of mental health traits using tweets.	222
177	The first model used Recurrent Neural Network	223
178	(RNN), while the remaining three used Attention,	224
179	which is the engine of transformer-based models.	225
180	We trained these model on the top of the PLM as	226
181	described in the section 2.3.	227
182	We used a single A100 80GB GPU to train all the	228
183	models. Each experiment took around 20 days	229
184	to complete. Hence, the limited number of experi-	
185	ments is due to the lack of resources for performing	
186	multiple experiments at once. Please note that the	
187	github link containing all the experiments (com-	
188	pleted and currently running) will be publicly avail-	
189	able in the final paper after the review process.	
190	3.1 Long Short Term Memory (LSTM)	
191	In our experiment, we implemented LSTM as our	
192	first experiment. Since the tweets are sequen-	
193	tial with each user having up to 1000 tweets and	
194	there are a differing number of tweets between the	
195	users, LSTM was appropriate as an experimental	
196	design. Further, LSTM stores long-term dependen-	
197	cies which fails on other neural networks (Hochre-	
198	iter and Schmidhuber, 1997). We implemented two	
199	LSTM models for this research with layers 1 and	
200	2. The added layer increased the complexity of	
201	the model. The number of hidden layers in both	
202	architectures were set to 100.	
203	3.2 Using attention mechanism	
204	Attention is the core of transformer based mod-	
205	els (Vaswani et al., 2017). Since we are using	
206	RoBERTa for the base model (Barbieri et al., 2020),	
207	which is a transformer based model, we added a	
208	multi-headed attention layer of 4 heads for our sec-	
	ond experiment design. This choice was made	209
	to attend to various parts of the tweet sequence	210
	differently. The idea behind this design was that	211
	the $[CLS]$ token would attend to a single tweet	212
	t_i and the stack of $[CLS]$ tokens from each user	213
	$t_n^{u_i}$ would use a cross-attention between the tweets	214
	i.e. $MHA(t_n^{u_i})$, where $MHA()$ is the multi-head	215
	attention. This would determine the presence or ab-	216
	sence of some mental health condition (depression	217
	or PTSD) for the user u_i .	218
	3.3 Two sentence sliding window	219
	For this experiment, we used two sentences ap-	220
	pended together before the tokenization i.e. for	221
	user u_i , $t_{u_i} = t_1 + t_2, t_2 + t_3, \dots, t_{n-1} + t_n$. A	222
	sliding window meant that except the first and the	223
	last tweet, every tweet in between would have in-	224
	formation linked with its previous and the next	225
	tweet, creating a short term attention. The result-	226
	ing $[CLS]$ token would go through cross-attention	227
	layer for long term attention across all the tweets	228
	belonging to a single user u_i , similar to section 3.2.	229
	3.4 Adding temporal information	230
	In this experiment, we added temporal information	231
	in terms of time lapse between the current and	232
	previous tweet as a part of the tweet. The first	233
	tweet t_1 was converted to $t_1 = "First\ tweet\ :"$	234
	$+ t_1$ and every subsequent tweets were converted	235
	to $t_i = "After\ x\ ;," + t_i$, where x was the time	236
	lapse between the current tweet t_i and the last tweet	237
	t_{i-1} , adding temporal context to the tweets. These	238
	were then processed in the same fashion as the	239
	attention as described in section 3.2.	240
	4 Evaluation	241
	We evaluated the experiments based on two key	242
	metrics: F1 score and Recall. The best performing	243
	results are presented in Table 1. Given p_1, p_2 and	244
	p_3 are probabilities for control, depression and	245
	PTSD respectively, the results were calculated as	246
	such for the mentioned three cases.	247
		248
	Case A. Multinomial classification: In this	249
	case, we performed the identification of control vs	250
	depression vs PTSD users based on the highest	251
	probability i.e. $max(p_1, p_2, p_3)$. This was the	252
	primary objective of this study.	253
		254
	Case B. Depression vs Control: In this case, we	255
	removed the probability p_3 from all experimental	256

Models	Multinomial (A)		Depression vs Control (B)		Mental Health vs Control (C)	
	F1 Score	Recall	F1 Score	Recall	F1 Score	Recall
LSTM (1 layer)	0.522	0.527	0.586	0.713	0.688	0.757
LSTM (2 layers)	0.504	0.524	0.569	0.827	0.712	0.913
Attention	0.595	0.590	0.616	0.780	0.724	0.830
Temporal	0.606	0.603	0.620	0.680	0.716	0.723
2 sentence Attention	0.635	0.637	0.655	0.760	0.755	0.797
MentalRoBERTa	-	-	0.697	0.703	-	-

Table 1: Performance metrics across experiments for control vs depression vs PTSD (multinomial) classification (A), depression vs control classification (B) and (depression or PTSD) vs control classification (C)

257 results and re-scaled the results for p_1 and p_2 and
258 evaluated using the readjusted probabilities. This
259 was done to compare our model results with the
260 baseline, MentalRoBERTa (Ji et al., 2021). Mental-
261 RoBERTa was taken as the baseline due to its large
262 scale training on mental health texts, including the
263 same data set as ours.

264 **Case C: Mental health vs Control:** In this case,
265 we added the probability of p_2 and p_3 from all
266 experimental results and evaluated using the new
267 probability. This was done to simulate a scenario
268 where presence or absence of any mental health
269 condition is tested. This also allowed us to confirm
270 or deny if there are overlapping sentiments among
271 users with depression and PTSD.

272 In our experiments, two-sentence attention
273 model achieved the best performance in both met-
274 rics for case A. Similarly, the same model per-
275 formed best in F1 score for case C, while recall
276 was higher for 2 layer LSTM for case C. Recall
277 was also higher for 2 layer LSTM in case B. How-
278 ever, for case B, our model did not outperform the
279 baseline F1 score of MentalRoBERTa model.

280 While our metrics are lower for case A in com-
281 parison to other cases, it is expected of a multinom-
282 inal classification compared to binary classification.
283 Identification of depression and PTSD separately
284 resulted in decreased performance, compared to
285 case B where only depression is identified and case
286 C where general mental health condition is identi-
287 fied. Another possible explanation is the potential
288 overlap of expressions in tweets from users with
289 depression and PTSD. Consequently, the classifi-
290 cation between the two groups becomes more chal-
291 lenging compared to the classification of an individ-
292 ual mental disorder from the control group alone.
293 However, when these disorders are combined, the
294 result improves significantly as seen from the re-

sults in case C of Table 1.

295 High values of recall in 2 layer LSTM for case
296 B (0.827) and case C (0.913) also means that ma-
297 jority of mental health users are identified. While
298 this causes less generalisation as demonstrated by
299 their corresponding F1 score, it is desirable for
300 this particular study because not identifying mental
301 health users are more costly than identifying false
302 positives of the same.

303 It should be noted that building state-of-the-art
304 model was not the primary objective of this study.
305 Instead, it was a study to target identification of
306 multiple mental health disorder for early diagnosis.
307 Further, these models cannot replace psychiatric
308 diagnosis and therapeutic interventions, but they
309 are valuable tools to aid clinicians and researchers.
310

311 5 Conclusion

312 In this study, we implemented and evaluated PLM
313 efficacy in identifying multiple mental health condi-
314 tions, including depression and PTSD from Twitter
315 data. Our experiments, including LSTM, atten-
316 tion, sliding window approach, and the integra-
317 tion of temporal information, showed that the two-
318 sentence attention model performs adequately for
319 detecting multiple health conditions. While the per-
320 formance was not as high as binary identification,
321 it can be attributed due to the overlap of sentiments
322 in tweets between depression and PTSD users. Our
323 findings also indicate that two layer LSTM model
324 is better at detecting the presence of depression
325 or mental health, in general, but it failed to gener-
326 alise well. In this regard, perhaps attention based
327 mechanism was significant as well. Despite lower
328 metrics in multi-nominal setting, our study pro-
329 vides an avenue of early mental health detection,
330 potentially leading to better targeted treatment and
331 interventions using social media.

332 Limitations

333 One of the aforementioned limitations is that only
334 last 1000 tweets (if more than 1000 tweets present)
335 per user were considered for this research. The
336 GPU server was shared between various projects as
337 well as the lack of resources to add more GPU
338 servers meant that not all tweets could be pro-
339 cessed. The reliance on the processing of tweets
340 sequentially further meant that each epoch was
341 much longer, since batching was not possible. This
342 caused each model to run up to 20 days, hence re-
343 sulting in lower number of experiments. Further,
344 only a single dataset was used, which could bias
345 the results. In addition, the tweets were extracted a
346 decade ago, which means the newer tweets would
347 not have been collected. The lexicon in which hu-
348 mans express sentiments perhaps changed in the
349 last decade and those were not captured. Addition-
350 ally, the collected tweets are only a sub-sample
351 of the much larger cohort of mental health users
352 who are not considered in this study. Even while
353 focusing on this cohort itself, there is a lack of evi-
354 dence to affirm the presence or absence of mental
355 health conditions between the Twitter users. Fi-
356 nally, our study aims to develop a model for as-
357 sisting researchers and clinicians for detection of
358 mental health conditions using social context for
359 non-clinical use. However, it does not replace clin-
360 ical diagnoses which is essential for the detection
361 and treatment of mental health issues.

362 Ethics Statement

363 The ethics was approved in accordance to Human
364 Research Ethics Committee (HREC) approval num-
365 ber H15559. The data was already de-identified
366 when it was received from Department of Com-
367 puter Science, John Hopkins University.

368 Acknowledgements

369 We would like to thank Professor Mark Dredze
370 for providing data and Department of Computer,
371 Data and Mathematical Sciences, Western Sydney
372 University for allowing to use the GPU clusters
373 which allowed for processing of data.

374 References

375 Francesco Barbieri, Jose Camacho-Collados, Leonardo
376 Neves, and Luis Espinosa-Anke. 2020. [TweetEval:
377 Unified Benchmark and Comparative Evaluation for
378 Tweet Classification](#). ArXiv:2010.12421 [cs].

- Diogo Beirão, Helena Monte, Marta Amaral, Alice Lon- 379
gras, Carla Matos, and Francisca Villas-Boas. 2020. 380
[Depression in adolescence: a review](#). *Middle East* 381
Current Psychiatry, 27(1):50. 382
- Janie Busby Grant, Philip J. Batterham, Sonia M. Mc- 383
Callum, Aliza Werner-Seidler, and Alison L. Calear. 384
2023. [Specific anxiety and depression symptoms](#) 385
[are risk factors for the onset of suicidal ideation and](#) 386
[suicide attempts in youth](#). *Journal of Affective Disor-* 387
ders, 327:299–305. 388
- Patricia A. Cavazos-Rehg, Melissa J. Krauss, Shaina 389
Sowles, Sarah Connolly, Carlos Rosas, Meghana 390
Bharadwaj, and Laura J. Bierut. 2016. [A content](#) 391
[analysis of depression-related tweets](#). *Computers in* 392
Human Behavior, 54:351–357. 393
- Po-Han Chou, Shao-Cheng Wang, Chi-Shin Wu, and 394
Masaya Ito. 2023. [Trauma-related guilt as a mediator](#) 395
[between post-traumatic stress disorder and suicidal](#) 396
[ideation](#). *Frontiers in Psychiatry*, 14. Publisher: 397
Frontiers. 398
- Glen Coppersmith, Mark Dredze, and Craig Harman. 399
2014. [Quantifying Mental Health Signals in Twitter](#). 400
In *Proceedings of the Workshop on Computational* 401
Linguistics and Clinical Psychology: From Linguistic 402
Signal to Clinical Reality, pages 51–60, Baltimore, 403
Maryland, USA. Association for Computational Lin- 404
guistics. 405
- Glen Coppersmith, Mark Dredze, Craig Harman, 406
Kristy Hollingshead, and Margaret Mitchell. 2015. 407
[CLPsych 2015 Shared Task: Depression and PTSD](#) 408
[on Twitter](#). In *Proceedings of the 2nd Workshop on* 409
Computational Linguistics and Clinical Psychology: 410
From Linguistic Signal to Clinical Reality, pages 31– 411
39, Denver, Colorado. Association for Computational 412
Linguistics. 413
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 414
Kristina Toutanova. 2018. Bert: Pre-training of deep 415
bidirectional transformers for language understand- 416
ing. *arXiv preprint arXiv:1810.04805*. Type: Journal 417
article. 418
- Son Doan, Amanda Ritchart, Nicholas Perry, Juan D 419
Chaparro, and Mike Conway. 2017. [How Do You](#) 420
[#relax When You’re #stressed? A Content Analysis](#) 421
[and Infodemiology Study of Stress-Related Tweets](#). 422
JMIR Public Health and Surveillance, 3(2):e35. 423
- Jon Finch. 2023. [The Difference Between Depression](#) 424
[and PTSD](#). 425
- Luis Roberto García-Noguez, Saúl Tovar-Arriaga, Wil- 426
frido Jacobo Paredes-García, Juan Manuel Ramos- 427
Arreguín, and Marco Antonio Aceves-Fernandez. 428
2023. [Automatic classification of depressive users](#) 429
[on Twitter including temporal analysis](#). *Network* 430
Modeling Analysis in Health Informatics and Bioin- 431
formatics, 12(1):38. 432
- S. Hochreiter and J. Schmidhuber. 1997. [Long short-](#) 433
[term memory](#). *Neural Computation*, 9(8):1735–1780. 434
Type: Journal article. 435

436 Ryan Holliday, Claire A. Hoffmire, W. Blake Martin, 492
437 Rani A. Hoff, and Lindsey L. Monteith. 2021. *As-* 493
438 *associations between justice involvement and PTSD* 494
439 *and depressive symptoms, suicidal ideation, and sui-*
440 *cide attempt among post-9/11 veterans.* *Psychologi-* 495
441 *cal Trauma: Theory, Research, Practice, and Policy,* 496
442 13(7):730–739. Place: US Publisher: Educational 497
443 Publishing Foundation. 498

444 Ahmed Husseini Orabi, Prasadith Buddhitha, Mahmoud 492
445 Husseini Orabi, and Diana Inkpen. 2018. *Deep* 493
446 *Learning for Depression Detection of Twitter Users.* 494
447 *In Proceedings of the Fifth Workshop on Computa-*
448 *tional Linguistics and Clinical Psychology: From* 495
449 *Keyboard to Clinic*, pages 88–97, New Orleans, LA. 496
450 Association for Computational Linguistics. 497

451 Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, 492
452 Prayag Tiwari, and Erik Cambria. 2021. *Mental-* 493
453 *BERT: Publicly Available Pretrained Language Mod-* 494
454 *els for Mental Healthcare.* ArXiv:2110.15621 [cs]. 495

455 Megan C. Kearns, Kerry J. Ressler, Doug Zatz- 492
456 ick, and Barbara Olasov Rothbaum. 2012. *Early* 493
457 *Interventions for Ptsd: A Review.* *De-* 494
458 *pression and Anxiety*, 29(10):833–842. _eprint: 495
459 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/da.21997>. 496

460 Kimberly Holland, Timothy J. Legg. 2019. *PTSD and* 492
461 *Depression: Similarities, Differences & What If You* 493
462 *Have Both.* 494

463 Layla Kratovic, Lia J. Smith, and Anka A. Vu- 492
464 janovic. 2021. *PTSD Symptoms, Suicidal* 493
465 *Ideation, and Suicide Risk in University Stu-* 494
466 *dents: The Role of Distress Tolerance.* *Jour-* 495
467 *nal of Aggression, Maltreatment & Trauma,* 496
468 30(1):82–100. Publisher: Routledge _eprint: 497
469 <https://doi.org/10.1080/10926771.2019.1709594>. 498

470 Andrew Page, Shu-Sen Chang, and David Gunnell. 492
471 2011. *Surveillance of Australian Suicidal Behaviour* 493
472 *Using the Internet?* *Australian & New Zealand* 494
473 *Journal of Psychiatry*, 45(12):1020–1022. Publisher: 495
474 SAGE Publications Ltd. 496

475 Vedant Vajre, Mitch Naylor, Uday Kamath, and Amarda 492
476 Shehu. 2021. *PsychBERT: A Mental Health Lan-* 493
477 *guage Model for Social Media Mental Health Be-* 494
478 *havioral Analysis.* In *2021 IEEE International Con-* 495
479 *ference on Bioinformatics and Biomedicine (BIBM),* 496
480 pages 1077–1082. 497

481 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob 492
482 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 493
483 Kaiser, and Illia Polosukhin. 2017. Attention is all 494
484 you need. *Advances in neural information processing* 495
485 *systems*, 30. Type: Journal article. 496

486 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien 492
487 Chaumond, Clement Delangue, Anthony Moi, Pier- 493
488 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, 494
489 et al. 2019. Huggingface’s transformers: State-of- 495
490 the-art natural language processing. *arXiv preprint* 496
491 *arXiv:1910.03771.* 497